

UNITED STATES DESIGN PATENT APPLICATION

FOR

**CHIPSET SUPPORT FOR MANAGING HARDWARE INTERRUPTS IN A VIRTUAL
MACHINE SYSTEM**

Inventors:

STALINSELVARAJ JEYASINGH**ANDREW V. ANDERSON****STEPHEN M. BENNETT****ERIK COTA-ROBLES****ALAIN KAGI****GILBERT NEIGER****RICHARD UHLIG**

Prepared by:

BLAKELY SOKOLOFF TAYLOR & ZAFMAN LLP
12400 Wilshire Boulevard, Seventh Floor
Los Angeles, CA 90025-1026

(408) 720-8300

EXPRESS MAIL CERTIFICATE OF MAILING"Express Mail" mailing label number EV341060355USDate of Deposit September 30, 2003

I hereby certify that this paper or fee is being deposited with the United States Postal Service "Express Mail Post Office to Addressee" service under 37 CFR 1.10 on the date indicated above and is addressed to the Mail Stop Patent Application, Commissioner for Patents, P.O. Box 1450, Alexandria, VA 22313-1450

Michelle Begay

(Typed or printed name of person mailing paper or fee)

Michelle B.

(Signature of person mailing paper or fee)

CHIPSET SUPPORT FOR MANAGING HARDWARE INTERRUPTS IN A VIRTUAL MACHINE SYSTEM

Field

[0001] Embodiments of the invention relate generally to virtual machines, and more specifically to managing hardware interrupts in a virtual machine system.

Background of the Invention

[0002] In a typical computer system, devices request services from system software by generating interrupt requests, which are propagated to an interrupt controller via multiple interrupt request lines. Once the interrupt controller identifies an active interrupt request line, it sends an interrupt signal to the processor. In response, the interrupt controller interface logic on the processor determines whether the software is ready to receive the interrupt. If the software is not ready to receive the interrupt, the interrupt is held in a pending state until the software becomes ready. Once the software is determined to be ready, the interrupt controller interface logic requests the interrupt controller to report which of the pending interrupts is highest priority. The interrupt controller prioritizes among the various interrupt request lines and identifies the highest priority interrupt request to the processor which then transfers control flow to the code that handles that interrupt request.

[0003] In a conventional operating system (OS), all the interrupts are controlled by a single entity known as an OS kernel. In a virtual machine system, a virtual-machine monitor (VMM) should have ultimate control over various operations and events occurring in the system to provide proper operation of virtual machines and for protection from and between virtual machines. To achieve this, the VMM typically receives control when guest software accesses certain hardware resources or certain events occur, such as an interrupt or an exception. In particular, when system devices generate interrupts, the VMM may intercede between the virtual machine and the interrupt controller device. That is, when an interrupt signal is raised, the currently running virtual machine is interrupted and control of the processor is passed to the VMM. The VMM then receives the interrupt and handles the interrupt or delivers the interrupt to an appropriate virtual machine.

Brief Description of the Drawings

[0004] The present invention is illustrated by way of example, and not by way of limitation, in the figures of the accompanying drawings and in which like reference numerals refer to similar elements and in which:

[0005] **Figure 1** illustrates one embodiment of a virtual-machine environment, in which the present invention may operate;

[0006] **Figure 2** is a block diagram of one embodiment of a system for processing interrupts in a virtual machine environment;

[0007] **Figure 3** is a flow diagram of one embodiment of a process for handling interrupts in a virtual machine system;

[0008] **Figure 4** is a flow diagram of one embodiment of a process for configuring the handling of interrupts during execution of a virtual machine in a virtual machine system; and

[0009] **Figure 5** is a flow diagram of one embodiment of a process for managing the configuration of the chipset interrupt control during a switch in virtual machines in a virtual machine system.

Description of Embodiments

[0010] A method and apparatus for controlling interrupts in a virtual machine system are described. In the following description, for purposes of explanation, numerous specific details are set forth in order to provide a thorough understanding of the present invention. It will be apparent, however, to one skilled in the art that the present invention can be practiced without these specific details.

[0011] Some portions of the detailed descriptions that follow are presented in terms of algorithms and symbolic representations of operations on data bits within a computer system's registers or memory. These algorithmic descriptions and representations are the means used by those skilled in the data processing arts to most effectively convey the substance of their work to others skilled in the art. An algorithm is here, and generally, conceived to be a self-consistent sequence of operations leading to a desired result. The operations are those requiring physical manipulations of physical quantities. Usually, though not necessarily, these quantities take the form of electrical or magnetic signals capable of being stored, transferred, combined, compared, and otherwise manipulated. It has proven convenient at times, principally for reasons of common usage, to refer to these signals as bits, values, elements, symbols, characters, terms, numbers, or the like.

[0012] It should be borne in mind, however, that all of these and similar terms are to be associated with the appropriate physical quantities and are merely convenient labels applied to these quantities. Unless specifically

stated otherwise as apparent from the following discussions, it is appreciated that throughout the present invention, discussions utilizing terms such as "processing" or "computing" or "calculating" or "determining" or the like, may refer to the action and processes of a computer system, or similar electronic computing device, that manipulates and transforms data represented as physical (electronic) quantities within the computer system's registers and memories into other data similarly represented as physical quantities within the computer-system memories or registers or other such information storage, transmission or display devices.

[0013] In the following detailed description of the embodiments, reference is made to the accompanying drawings that show, by way of illustration, specific embodiments in which the invention may be practiced. In the drawings, like numerals describe substantially similar components throughout the several views. These embodiments are described in sufficient detail to enable those skilled in the art to practice the invention. Other embodiments may be utilized and structural, logical, and electrical changes may be made without departing from the scope of the present invention. Moreover, it is to be understood that the various embodiments of the invention, although different, are not necessarily mutually exclusive. For example, a particular feature, structure, or characteristic described in one embodiment may be included within other embodiments. The following detailed description is, therefore, not to be taken in a limiting sense, and the

scope of the present invention is defined only by the appended claims, along with the full scope of equivalents to which such claims are entitled.

[0014] Although the below examples may describe embodiments of the present invention in the context of execution units and logic circuits, other embodiments of the present invention can be accomplished by way of software. For example, in some embodiments, the present invention may be provided as a computer program product or software which may include a machine or computer-readable medium having stored thereon instructions which may be used to program a computer (or other electronic devices) to perform a process according to the present invention. In other embodiments, steps of the present invention might be performed by specific hardware components that contain hardwired logic for performing the steps, or by any combination of programmed computer components and custom hardware components.

[0015] Thus, a machine-readable medium may include any mechanism for storing or transmitting information in a form readable by a machine (e.g., a computer), but is not limited to, floppy diskettes, optical disks, Compact Disc, Read-Only Memory (CD-ROMs), and magneto-optical disks, Read-Only Memory (ROMs), Random Access Memory (RAM), Erasable Programmable Read-Only Memory (EPROM), Electrically Erasable Programmable Read-Only Memory (EEPROM), magnetic or optical cards, flash memory, a transmission over the Internet, electrical, optical, acoustical or other forms of propagated signals (e.g., carrier waves, infrared signals, digital signals, etc.) or the like.

[0016] Further, a design may go through various stages, from creation to simulation to fabrication. Data representing a design may represent the design in a number of manners. First, as is useful in simulations, the hardware may be represented using a hardware description language or another functional description language. Additionally, a circuit level model with logic and/or transistor gates may be produced at some stages of the design process. Furthermore, most designs, at some stage, reach a level of data representing the physical placement of various devices in the hardware model. In the case where conventional semiconductor fabrication techniques are used, data representing a hardware model may be the data specifying the presence or absence of various features on different mask layers for masks used to produce the integrated circuit. In any representation of the design, the data may be stored in any form of a machine-readable medium. An optical or electrical wave modulated or otherwise generated to transmit such information, a memory, or a magnetic or optical storage such as a disc may be the machine readable medium. Any of these mediums may "carry" or "indicate" the design or software information. When an electrical carrier wave indicating or carrying the code or design is transmitted, to the extent that copying, buffering, or re-transmission of the electrical signal is performed, a new copy is made. Thus, a communication provider or a network provider may make copies of an article (a carrier wave) embodying techniques of the present invention.

[0017] **Figure 1** illustrates one embodiment of a virtual-machine environment 100, in which the present invention may operate. In this embodiment, bare platform hardware 116 comprises a computing platform, which may be capable, for example, of executing a standard operating system (OS) or a virtual-machine monitor (VMM), such as a VMM 112.

[0018] The VMM 112, though typically implemented in software, may emulate and export a bare machine interface to higher level software. Such higher level software may comprise a standard or real-time OS, may be a highly stripped down operating environment with limited operating system functionality, may not include traditional OS facilities, etc. Alternatively, for example, the VMM 112 may be run within, or on top of, another VMM. VMMs may be implemented, for example, in hardware, software, firmware or by a combination of various techniques.

[0019] The platform hardware 116 can be of a personal computer (PC), mainframe, handheld device, portable computer, set-top box, or any other computing system. The platform hardware 116 includes a processor 118, memory 120, chipset core logic 122 and one or more interrupt sources 128.

[0020] Processor 118 can be any type of processor capable of executing software, such as a microprocessor, digital signal processor, microcontroller, or the like. The processor 118 may include microcode, programmable logic or hardcoded logic for performing the execution of method embodiments of the present invention. Though **Figure 1** shows only one such processor 118, there may be one or more processors in the system.

[0021] Memory 120 can be a hard disk, a floppy disk, random access memory (RAM), read only memory (ROM), flash memory, any combination of the above devices, or any other type of machine medium readable by processor 118. Memory 120 may store instructions and/or data for performing the execution of method embodiments of the present invention.

[0022] The one or more interrupt sources 128 may be, for example, input-output (I/O) devices (e.g., network interface cards, communication ports, video controllers, disk controllers) on system buses (e.g., PCI, ISA, AGP) or devices integrated into the chipset logic or processor (e.g., real-time clocks, programmable timers, performance counters).

[0023] The VMM 112 presents to other software (i.e., "guest" software) the abstraction of one or more virtual machines (VMs), which may provide the same or different abstractions to the various guests. **Figure 1** shows two VMs, 102 and 114. The guest software running on each VM may include a guest OS such as a guest OS 104 or 106 and various guest software applications 108 and 110. Each of the guest OSs 104 and 106 expect to access physical resources (e.g., processor registers, memory and I/O devices) within the VMs 102 and 114 on which the guest OS 104 or 106 is running and to handle various events including interrupts generated by system devices during the operation of the VMs 102 and 114.

[0024] Some interrupts may need to be handled by a currently operating VM. Other interrupts may need to be handled by the VMM 112 or a VM that is not currently operating. If the interrupt is to be handled by the

currently-operating VM, control remains with this VM, and the interrupt is delivered to this VM if it is ready to receive interrupts (as indicated, for example, by an interrupt flag in a designated processor register). If the interrupt is to be handled by the VMM 112, control is transferred to the VMM 112. The transfer of control from guest software to the VMM 112 is referred to herein as a VM exit. After receiving control following the VM exit, the VMM 112 may perform a variety of processing, including, for example, acknowledging and handling the interrupt, after which it may return control to guest software. If the VMM does not handle the interrupt itself, it may facilitate delivery of the interrupt to a VM designated to handle the interrupt. The transfer of control from the VMM to guest software is referred to as a VM entry.

[0025] In one embodiment, the processor 118 controls the operation of the VMs 102 and 114 in accordance with data stored in a virtual machine control structure (VMCS) 124. The VMCS 124 is a structure that may contain state of guest software, state of the VMM 112, execution control information indicating how the VMM 112 wishes to limit or otherwise control operation of guest software, information controlling transitions between the VMM 112 and a VM, etc. In one embodiment, the VMCS 124 is stored in memory 120. In another embodiment, the VMCS 124 is stored in the processor 118. In some embodiments, multiple VMCS structures are used to support multiple VMs.

[0026] The processor 118 reads information from the VMCS 124 to determine the execution environment of the VM and to constrain its behavior.

For example, the processor 118 may consult the execution control information in the VMCS to determine if external interrupts are to cause VM exits. When a VM exit occurs, components of the processor state used by guest software are saved to the VMCS 124, and components of the processor state required by the VMM 112 are loaded from the VMCS 124. When a VM entry occurs, the processor state that was saved at the VM exit is restored using data stored in the VMCS 124, and control is returned to guest software.

[0027] In one embodiment, chipset core logic 122 is coupled to the processor 118 to assist the processor 118 in handling interrupts generated by the interrupt sources 128, described below. The chipset logic 122 may be part of the processor 118 or an independent component. The chipset logic 122 may include microcode, programmable logic or hardcoded logic for performing the execution of method embodiments of the present invention. Though Figure 1 shows only one chipset logic 122, the system may contain one or more such chipsets.

[0028] As will be discussed in more detail below, the chipset core logic 122 controls the distribution of interrupts generated by the one or more interrupt sources 128. If an interrupt is generated by a device that is managed by a currently operating VM, the chipset core logic 122 sends the interrupt request to the processor 118 for delivery to the currently operating VM. If an interrupt is generated by a device that is not managed by a currently operating VM, the chipset core logic 122 either holds the interrupt pending or

indicates to the processor 118 that control is to be transitioned to the VMM 112.

[0029] Figure 2 is a block diagram of one embodiment of a system 200 for processing interrupts in a virtual machine environment. The system 200 includes a processor 220 and chipset core logic 202. The system 200 is active during operation of the VMM and during operation of guest software. The term “currently operating software” is used herein to refer to the VMM or the currently operating guest software (i.e., the currently operating VM).

[0030] In one embodiment, the chipset core logic 202 includes an interrupt controller 204, a set of multiplex blocks 206, and a VMM block 210.

[0031] The interrupt controller 204 receives interrupt request signals generated by system devices and delivers interrupt requests (INTRs) to the processor 220. In an embodiment, the processor acknowledges an interrupt request from the interrupt controller 204 by an interrupt acknowledgement bus cycle (INTA), for which the interrupt controller 204 returns the vector number of the interrupt service routine to be executed to service the interrupt.

[0032] The interrupt controller 204 has a state machine and a number of registers, which may include readable and writable registers, read-only registers and write-only registers. In one embodiment, the interrupt controller 204 includes a read and write access path to its registers to allow a VMM to save the current state of the interrupt controller 204 (e.g., in the VMCS or in a VMM data structure) and to restore the interrupt controller state associated with a VM that is to be invoked. In addition, in one embodiment, the current

state of the state machine in the interrupt controller 204 is readable and writable to allow the VMM to switch VMs. In another embodiment, the state machine in the interrupt controller 204 is readable and the VMM has to replay a series of writes to the interrupt controller 204 to place the interrupt controller 204 into the correct state.

[0033] The interrupt controller 204 may be a conventional interrupt controller (e.g., an 8259A interrupt controller) that is enhanced to provide a read and write access path to its registers and state machine. The read and write access path may be implemented, for example, as an extension to chipset specific registers to provide the behavior that is identical to the behavior of the conventional interrupt controller (e.g., the read-only registers remain read-only unless being accessed by the VMM). This can be achieved, for example, by providing access to the registers of the interrupt controller through memory mapped I/O registers. The interrupt controller 204 may include masking and prioritization logic that is known to one of ordinary skill in the art.

[0034] The interrupt controller 204 is coupled to the multiplex blocks 206. The number of multiplex blocks 206 is equal to the number of interrupt request lines 208. Each interrupt request line 208 propagates interrupt requests generated by a specific device to a corresponding multiplex block 206.

[0035] In a virtual machine system, a device may be managed by a certain VM and generate interrupts that need to be delivered to this VM. For

example, a video capture card may be managed directly by a single VM and may not be visible to other VMs. Alternatively, a device may be managed by several VMs and generate interrupts that need to be delivered to multiple VMs. Yet, alternatively, the device may be managed by the VMM and generate interrupts that need to be delivered to the VMM. For example, an integrated drive electronics (IDE) controller may be managed exclusively by the VMM.

[0036] Each multiplex block 206 receives interrupt request signals generated by a specific device. Depending on its setting, each multiplex block 206 may route an interrupt request signal to the interrupt controller 204 or the VMM block 210. In one embodiment, the multiplex blocks 206 are set by the VMM. The VMM may set the multiplex blocks 206 prior to requesting a VM entry and/or following a VM exit. In another embodiment, the multiplex blocks 206 are set by the processor as part of a VM entry and/or a VM exit. The values used to set the multiplex blocks 206 may be stored in the VMCS or in any other data structure.

[0037] In one embodiment illustrated in **Figure 2**, when a multiplex block 206 is set to route interrupt request signals generated by a corresponding device to the interrupt controller 204, the interrupt request signal may be routed to a specific input of the interrupt controller 204. For example, **Figure 2** shows interrupt IRQ0 being routed to input 1 of the interrupt controller 204 and IRQ1 is shown routed to input 0 of the interrupt

controller 204. Other embodiments allow the interrupt request signal to be routed to a single fixed input of the interrupt controller 204.

[0038] In addition, in the embodiment illustrated in **Figure 2**, when a multiplex block 206 is set to route interrupt request signals generated by a corresponding device to the VMM block 210, the interrupt request signal can be routed to a fixed input of the VMM block 210. For example, in the example shown in **Figure 2**, IRQ_n is routed to the nth input of the VMM block 210. In another embodiment, the interrupt request signal may be routed by the multiplex block 206 to an arbitrary input of the VMM block 210.

[0039] The VMM block 210 includes a number of input lines connected to the output lines of the multiplex blocks 206. The VMM block 210 generates an output signal indicating to the processor 200 that control needs to be transitioned to the VMM. In one embodiment, the output signal is generated by combining the interrupt request signals routed to the VMM block 210 with an external signal generated by an external signal source 218. An external signal may be any signal, other than a signal of an external interrupt type, that may be configured by the VMM to cause a VM exist. For example, as shown in **Figure 2**, an external signal may be a non-maskable interrupt (NMI) signal that is configured to cause a VM exit each time the NMI occurs. The interrupt request signals are combined with the external signal using a Boolean OR operator illustrated by a gate 216.

[0040] In one embodiment, the VMM block 210 includes a mask register 212 that can mask an interrupt request signal routed to the VMM block 210.

In one embodiment, the VMM configures the mask register 212 to allow masking of interrupt request signals generated by certain devices. Masking may be used when the VMM does not need to be notified about interrupts generated by a device managed exclusively by a non-currently operating VM. In one embodiment, the mask register 212 is set by the VMM. The VMM may set the mask register 212 prior to requesting a VM entry and/or following a VM exit. In another embodiment, the mask register 212 is set by the processor as part of a VM entry and/or a VM exit. The value used to set the mask register 212 may be stored in the VMCS or in any other data structure.

[0041] In one embodiment, the VMM block 210 includes a status register 214 that stores status of interrupt input lines of the VMM block 210. In one embodiment, when an interrupt request signal is asserted on an interrupt input line of the VMM block 210, a bit associated with this interrupt input line in the status register 214 is set. The status register 214 may be read by the VMM to obtain the status of the interrupt input lines of the VMM block 210. The status may be used, for example, to determine whether the output signal sent to the processor 220 by the VMM block 210 resulted from an interrupt propagated from an interrupt request line 208 or from an external signal (e.g., an NMI signal) generated by the external signal source 218.

[0042] As discussed above, in one embodiment, the VMM configures the multiplex blocks 206 prior to requesting a transfer of control to a VM. In this embodiment, once the control is transferred to a VM as requested by the VMM, only interrupt request signals generated by devices managed by the

VM are allowed to reach the interrupt controller 204. Interrupts managed by the VMM are routed to the VMM block 210. Interrupts managed by other VMs are routed to the VMM block and may be masked using masking register 212. Thus, the chipset core logic 202 allows a currently-operating VM to control all the interrupts generated by the devices managed by the currently-operating VM, while allowing the VMM to gain control over interrupts generated by devices owned by the VMM and other VMs.

[0043] In some embodiments, the chipset core logic 202 may include multiple interrupt controllers 204, with each interrupt controller 204 serving a corresponding VM. In those embodiments, the VMM does not need to save and restore the state of the interrupt controller each time it switches from one VM to another.

[0044] As discussed above, in some embodiments, the multiplex blocks 206 and mask register 212 are configured as part of a VM entry (e.g., by the VMM before requesting a VM entry or by the processor when performing a VM entry). Hence, the settings of the multiplex blocks 206 and mask register 212 remain unchanged following a VM exit, and interrupts generated during the operation of the VMM are routed according to the settings of a previously-operating VM. That is, if interrupts are not masked by the VMM, they may be delivered to the VMM via the interrupt controller 204 or the VMM block 210 using the mechanisms discussed above. The designated software within the VMM then handles these interrupts appropriately.

[0045] In other embodiments, the multiplex blocks 206 and mask register 212 are configured as part of a VM exit (e.g., by the VMM following a VM exit or by the processor when performing a VM exit). Hence, the multiplex blocks 206 and mask register 212 may be changed to route interrupts generated during the operation of the VMM differently than during operation of the VM which was running previously.

[0046] **Figure 3** is a flow diagram of one embodiment of a process 300 for handling interrupts in a virtual machine system. The process may be performed by processing logic that may comprise hardware (e.g., circuitry, dedicated logic, programmable logic, microcode, etc.), software (such as run on a general purpose computer system or a dedicated machine), or a combination of both. In one embodiment, processing logic is implemented in chipset core logic 202 of **Figure 2**.

[0047] Referring to **Figure 3**, process 300 begins with processing logic receiving, at a multiplex block, an interrupt request signal generated by a device (processing block 302). As discussed above, each multiplex block is coupled with a distinct interrupt request line that delivers interrupt request signals generated by a certain device.

[0048] At decision box 304, processing logic determines whether the interrupt request signal is to be routed to the interrupt controller. In one embodiment, this determination is done based on the configuration performed by the VMM.

[0049] If the determination made at decision box 304 is positive, processing logic sends the interrupt request signal to the interrupt controller (processing block 306), which may, according to the architecture of the interrupt controller, send the interrupt request to the processor (processing block 308). The interrupt controller may perform masking and prioritization of interrupts.

[0050] If the determination made at decision box 304 is negative, processing logic routes the interrupt request signal to the VMM block (processing block 310) and updates a status register of the VMM block to indicate that a signal has been asserted on an interrupt input line of the VMM block (processing block 312). Next, processing logic determines whether the interrupt request signal routed to the VMM block is masked (decision box 314). In one embodiment, the determination is made based on data set by the VMM in a mask register of the VMM block. If the interrupt request signal is masked, processing logic holds the interrupt pending (processing block 316). The interrupt may be held pending, for example, until a VM managing the device that generated this interrupt is invoked. At the time the VM managing the device is invoked, the VMM will modify the configuration of the multiplex blocks to route the interrupt request to the interrupt controller and hence allow the VM to manage the device directly.

[0051] If the interrupt request signal is not masked, processing logic generates an internal signal (processing block 318), combines the internal signal with an external signal generated by a designated external signal source

(e.g., an NMI source) using the OR operator (processing block 320), and delivers the resulting output signal to the processor (processing block 322).

[0052] The output signal may cause the processor to transfer control to the VMM. The VMM may then read the status register of the VMM block, determine that the signal resulted from an external interrupt, and handle this external interrupt itself or invoke an appropriate VM to handle it. It should be noted that an external signal (e.g., NMI) causing the output signal may require some predefined operations to be performed by the processor when asserted. If these predefined operations are different from the desired behavior, then the processor blocks the output signal to avoid the occurrence of the predefined operations. For example, an NMI may cause a transition within the VMM (e.g., by vectoring the NMI) when it is asserted during operation of the VMM. Blocking the NMI signal following the VM exit caused by the assertion of the NMI prevents the occurrence of the predefined transition within the VMM.

[0053] **Figure 4** is a flow diagram of one embodiment of a process 400 for handling interrupts in a virtual machine system. The process may be performed by processing logic that may comprise hardware (e.g., circuitry, dedicated logic, programmable logic, microcode, etc.), software (such as run on a general purpose computer system or a dedicated machine), or a combination of both. In one embodiment, processing logic is implemented in a VMM such as the VMM 112 of **Figure 1**.

[0054] Referring to **Figure 4**, process 400 begins with processing logic identifying interrupt request lines that are coupled to devices managed by a VM to be invoked (processing block 401). In one embodiment, these interrupt request lines are identified using data stored in the VMCS or any other data structure.

[0055] Next, processing logic configures a set of multiplex blocks (processing block 402). Based on this configuration, only interrupt request signals from the devices managed by the VM to be invoked can reach the interrupt controller. All the remaining interrupt signals are to be routed to the VMM block. In addition, in another embodiment, the VMM may configure the multiplex blocks to route the interrupt to a particular input of the VMM block.

[0056] At processing block 404, processing logic configures a mask register in the VMM block to allow selective masking of interrupt request signals routed to the VMM block. The selective masking can be used, for example, to hold some interrupt requests pending. Interrupt requests may be held pending if, for example, the VMM does not need to be notified about these interrupt requests because they come from a device that is managed exclusively by another VM. However, some interrupts belonging to another VM may not be masked if, for example, the other VM needs to run at a higher priority than the currently running VM. An example of such a situation is a VM which maintains some real-time quality of service with respect to some device (e.g., all interrupts need to be serviced in a specific amount of time).

[0057] At processing block 408, processing logic sets designated execution control fields in the VMCS to allow the VM being invoked to control I/O accesses to the interrupt controller and delivery of hardware interrupts. For example, the VMM may set designated execution control fields in the VMCS such that accesses to the I/O ports of the interrupt controller do not cause VM exits.

[0058] At processing block 410, processing logic sets a designated execution control field in the VMCS to cause a VM exit on each event generated when the VMM block asserts its output signal to the CPU. In one embodiment, the event may be a non-maskable interrupt (NMI). Other embodiments may use other events, according to the architecture of the chipset core logic and system. The VMM may also need to set controls such that certain predefined operations following VM exits are prevented from happening. For example, if the output signal of the VMM block of **Figure 2** is coupled to the NMI input of the processor, then the NMI signal may need to be blocked following VM exits that result from the assertion of the NMI signal.

[0059] At processing block 412, processing logic executes a VM entry instruction to request a transfer of control to the VM. Any mechanism in the art may be used to facilitate this transfer of control to the VM. After the transfer of control to the VM, the VM can directly access the interrupt controller (e.g., using I/O operations) and interrupts routed through the

interrupt controller will be handled directly by the VM, with no intervention by the VMM.

[0060] **Figure 5** is a flow diagram of one embodiment of a process 500 for managing chipset core logic state during a switch from one VM to another VM. The process may be performed by processing logic that may comprise hardware (e.g., circuitry, dedicated logic, programmable logic, microcode, etc.), software (such as run on a general purpose computer system or a dedicated machine), or a combination of both. In one embodiment, processing logic is implemented in a VMM such as the VMM 112 of **Figure 1**.

[0061] Referring to **Figure 5**, process 500 begins with processing logic recognizing that a transfer of control from one VM to another VM is pending (processing block 510). This may occur, for example, if the VMM is managing more than one VM, providing time slices to each VM in turn, similarly to how a traditional operating system may time slice processes on a single CPU.

[0062] At processing block 520, processing logic saves the current state of the interrupt controller. The state may be saved in the VMCS or any other designated data structure. In an embodiment, this saving of state is performed by the VMM. In another embodiment, the saving of the interrupt controller state is performed by the processor. In one embodiment, the savings of the interrupt controller state is performed as part of the processing at the time of a VM exit. Additionally, processing logic may need to save the state of the multiplex controls, and the VMM block masking register. The

saving of state relies on having the ability to read all appropriate state in the chipset core logic, as described above.

[0063] At processing block 540, processing logic restores the chipset core logic state from the previous operation of the VM to be activated. This restoration includes writing all appropriate state to the interrupt controller to configure the masking and control registers and state machine configuration that was present when the VM to be activated was last active. Additionally, the multiplex blocks and VMM block controls (e.g., masking register) in the chipset core logic may need to be reconfigured for the new VM, as described above. This restoration of state relies on having the ability to write to all appropriate state in the chipset core logic. In one embodiment, this restoration of chipset core logic state is performed by the VMM. In another embodiment, it is performed by the processor as part of the VM entry to the new VM.

[0064] At processing block 560, the VMM requests a transfer of control to the VM. In one embodiment, the VMM executes an instruction to initiate the transfer. Any mechanism in the art may be used to facilitate this transfer of control to the VM.

[0065] It should be noted that process 500 assumes that the execution controls were set appropriately prior to the first entry to the new VM (as described with regard to process 400).

[0066] Thus, a method and apparatus for handling interrupts in a virtual machine system have been described. It is to be understood that the

above description is intended to be illustrative, and not restrictive. Many other embodiments will be apparent to those of skill in the art upon reading and understanding the above description. The scope of the invention should, therefore, be determined with reference to the appended claims, along with the full scope of equivalents to which such claims are entitled.